

IN THE CLAIMS

Please amend the claims as follows:

1. (currently amended) A method for browser-enhanced web crawling associated with a network of hub processing units coupled to a plurality of information processing units over a network, the method executed by a web crawler on a hub processing unit associated with the network comprising the steps of:

retrieving a web document at an address, and extracting contents of the web document for rendering an intermediate dynamically constructed in-memory webpage representation of the web document at a hub processing unit which is formatted as if displayed for viewing on an end-user's web browser;

loading secondary documents associated with the web document in order to render the secondary documents as part of the in-memory webpage representation, wherein the secondary documents include one or more images with textual content embedded therein;

~~sending to one or more information processing units a browser side script to gather metadata; and performing the sub-steps of:~~

~~producing a final HTML markup;~~

analyzing and summarizing the in-memory webpage representation to produce a text map for the web page document of the textual contents therein~~final HTML markup to produce metadata; and~~

using optical character recognition on the images to extract textual content for adding to the textual map for the webpage document.

2. (currently amended) The method as defined in claim 1, wherein the retrieving ~~at the~~ web document at an address ~~step~~ further comprises retrieving a document at an address selected from the group of addresses consisting of a nodal address, a network address, a URL and equivalents.

3. (currently amended) The method as defined in claim 1, wherein the ~~analyzing and summarizing step further comprises analyzing and summarizing the whole and complete document~~ one or more images with textual content embedded therein include at least one of an in-line GIF image and an in-line JPEG image.

4. (currently amended) The method as defined in claim 1, ~~further comprising the step of analyzing any image data present in the document and any image data present in the documents utilizing optical character recognition techniques wherein the loading~~ secondary documents further comprises the loading of secondary documents including one or more Java applets with textual content embedded therein.

5. (currently amended) The method as defined in claim 1, wherein the ~~step of loading~~ secondary documents further comprises the loading of secondary documents including web documents selected from the group of documents consisting of in-line frames, frames, ~~and~~, images, image maps, applets, audio, video or equivalents.

6. (currently amended) The method as defined in claim 4, wherein the ~~step of analyzing any image data present in the document and any image data present in the documents utilizing optical character recognition techniques further comprises analyzing any images and image maps in the image data to produce text data~~ loading secondary documents further comprises the loading of secondary documents including one or more Java Script components with textual content embedded therein.

7. (currently amended) The method as defined in claim 1, wherein the retrieving step the web document further comprises performing the following sub-steps of:

initializing a first list with seed values;

checking if there are any URLs to be processed and ~~if there are, in response that any URL exists to be processed then~~ performing the following secondary-sub-steps of:

determining if a URL is in a second list; and ~~in response that a URL is not in the second list if it is not in the second list;~~ then performing the following tertiary-sub-steps of:

inserting the URL into the first list;

scheduling the URL for crawling;

crawling the URL when scheduled to do so;

removing the URL from the first list after the scheduled crawling;

entering the URL into the second list; and

repeating the checking step until there are no more URLs to be processed;

where if the determining step determines that the URL is in the second list then repeating the checking step until there are no more URLs to be processed.

8. (original) The method as defined in claim 7, wherein the sub-step of initializing a first list with seed values further includes the list being a URL pool.

9. (original) The method as defined in claim 7, wherein the sub-step of determining if a URL is in a second list further includes the second list being a visited pool.

10. (currently amended) The method as defined in claim 7, wherein the ~~tertiary~~-sub-step of crawling further comprises the sub-steps of:

issuing an HTTP command to a web server named in the URL;

receiving contents of an HTML page as a result of the issued HTTP command;

and

passing on the contents of the HTML page to a Page Rendering subroutine.

11. (original) The method as defined in claim 10, further including the sub-steps performed by the Page Rendering subroutine comprising:

- receiving the contents of the HTML page in the Page Rendering subroutine;
- building an in-memory representation of a Layout for the HTML page and if more data is needed to properly form the representation, then performing the sub-steps of:
 - requesting additional web-based information;
 - gathering this additional web-based information;
 - inserting any URLs associated with this additional web-based information into the second list and a URL cache;
 - building a final amended representation; and
 - forwarding the final amended representation to an Extraction subroutine;

wherein, if no more data is needed to properly form the in-memory representation, then forwarding the in-memory representation to the Extraction subroutine.

12. (original) The method as defined in claim 11, further including the sub-steps performed by the Page Extraction subroutine comprising:

- accessing a set of memory structures of the Page Renderer;
- copying a text portion of the structures into a text map;
- inspecting any in-line GIF and JPEG image references in the memory structures;
- extracting alternate text attributes;
- adding the alternate text attributes to a text map;
- invoking an optical character recognition engine;
- analyzing any in-line GIF and JPEG images using the optical character recognition engine for text content;
- extracting text content from the GIF and JPEG images;
- adding text content from the images to the text map; and
- forwarding the text map to a Page Summarizer subroutine.

13. (original) The method as defined in claim 12, further including the sub-steps performed by the Page Summarizer subroutine comprising:

ARC9-2000-0046-US1

5

09/607,370

receiving a text map from the Page Extractor subroutine;
processing the text map in an application-specific manner;
applying data extraction patterns to the text map;
translating resultant data from the applying step;
forwarding any URLs present in the text map to a manager subroutine; and
forwarding any extracted data and metadata to application logic.

14. (currently amended) A computer readable medium including programming instructions, the programming instructions including instructions for browser-enhanced web crawling associated with a network of hub processing units coupled to a plurality of information processing units over a network, the browser enhanced web crawling instructions on the computer readable medium comprising:

retrieving instructions for retrieving a web document at an address, and extracting contents of the web document for rendering an intermediate dynamically constructed in-memory webpage representation of the web document at a hub processing unit which is formatted as if displayed for viewing on an end-user's web browser;

loading instructions for loading secondary documents associated with the web document in order to render the secondary documents as part of the in-memory webpage representation, wherein the secondary documents include one or more images with textual content embedded therein;

sending instructions for sending to one or more information processing units a browser-side script to gather metadata;

producing instructions for producing a final HTML markup;
analyzing and summarizing instructions for analyzing and summarizing the in-memory webpage representation to produce a text map for the web page document of the textual contents therein final HTML markup to produce the final metadata; and
using optical character recognition on the images to extract textual content for adding to the textual map for the webpage document.

ARC9-2000-0046-US1

6

09/607,370

15. (currently amended) The computer readable medium as defined in claim 14, wherein the retrieving instructions for retrieving a web document at an address further comprises retrieving instructions for retrieving a document at an address selected from the group of addresses consisting of a nodal address, a network address, a URL and equivalents.

16. (currently amended) The computer readable medium as defined in claim 14, wherein the ~~analyzing and summarizing instructions further comprise analyzing and summarizing instructions for analyzing and summarizing the whole and complete document~~one or more images with textual content embedded therein include at least one of an in-line GIF image and an in-line JPEG image.

17. (currently amended) The computer readable medium as defined in claim 14, ~~further comprising image analyzing instructions for analyzing any image data present in the document and any image data present in the documents utilizing optical character recognition techniques~~wherein the loading secondary documents further comprises the loading of secondary documents including one or more Java applets with textual content embedded therein.

18. (currently amended) The computer readable medium as defined in claim 14, wherein the loading instructions for loading secondary documents further comprises loading instructions for loading of secondary documents including web documents selected from the group of documents consisting of in-line frames, frames, and images, ~~image maps, applets, audio, video or equivalents.~~

19. (currently amended) The computer readable medium as defined in claim 17, wherein the ~~analyzing instructions for analyzing any image data present in the document and any image data present in the documents utilizing optical character recognition techniques~~further comprises analyzing instructions for analyzing any images and image maps in the image data to produce text dataloading secondary

ARC9-2000-0046-US1

7

09/607,370

documents further comprises the loading of secondary documents including one or more Java Script components with textual content embedded therein.

20. (currently amended) A browser-enhanced web crawling unit associated with a network of a plurality of hub processing units coupled to a plurality of information processing units over a network, the browser enhanced web crawling unit on a hub processing unit comprising:

a retrieval unit for retrieving a web document at an address, and extracting contents of the web document for rendering an intermediate dynamically constructed in-memory webpage representation of the web document at a hub processing unit which is formatted as if displayed for viewing on an end-user's web browser;

a loader for loading secondary documents as required associated with the web document in order to render the secondary documents as part of the in-memory webpage representation, wherein the secondary documents include one or more images with textual content embedded therein;

~~an output for sending to one or more information processing units a browser side script to gather metadata;~~

~~a producer for producing a final HTML markup; and~~

a summarizer for analyzing and summarizing the in-memory webpage representation to produce a text map for the web page document of the textual contents therein~~final HTML markup to produce the final metadata~~

~~a producer for producing a final HTML markup;~~

an optical character recognition engine for use on the images to extract textual content for adding to the textual map for the webpage document.